

# Hybrid Machine Learning and Ensemble Data Assimilation

Wei Pan

**Co-authors and contributors:** Massimo Bonavita, Elías Hólm, Marcin Chrust, Alban Farchi, Marc Bocquet, Patrick Laloyaux, Ivo Pasmans, Sebastien Massart, Boštjan Melinc, Žiga Zaplotnik, Simon Lang, Ewan Pinnington

5<sup>th</sup> ECMWF-ESA ML Workshop, Bologna, 13-17 April 2026

# Ensemble 4D-Var

4DVar in Bayes form:

$$p(\mathbf{x} | \mathbf{y}_{\text{obs}}) \propto \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}_b)\right) p(\mathbf{y}_{\text{obs}} | \mathbf{x})$$

Variational DA produces *maximum a posteriori* (MAP) estimates

$$\mathbf{x}^a = \arg \max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}).$$

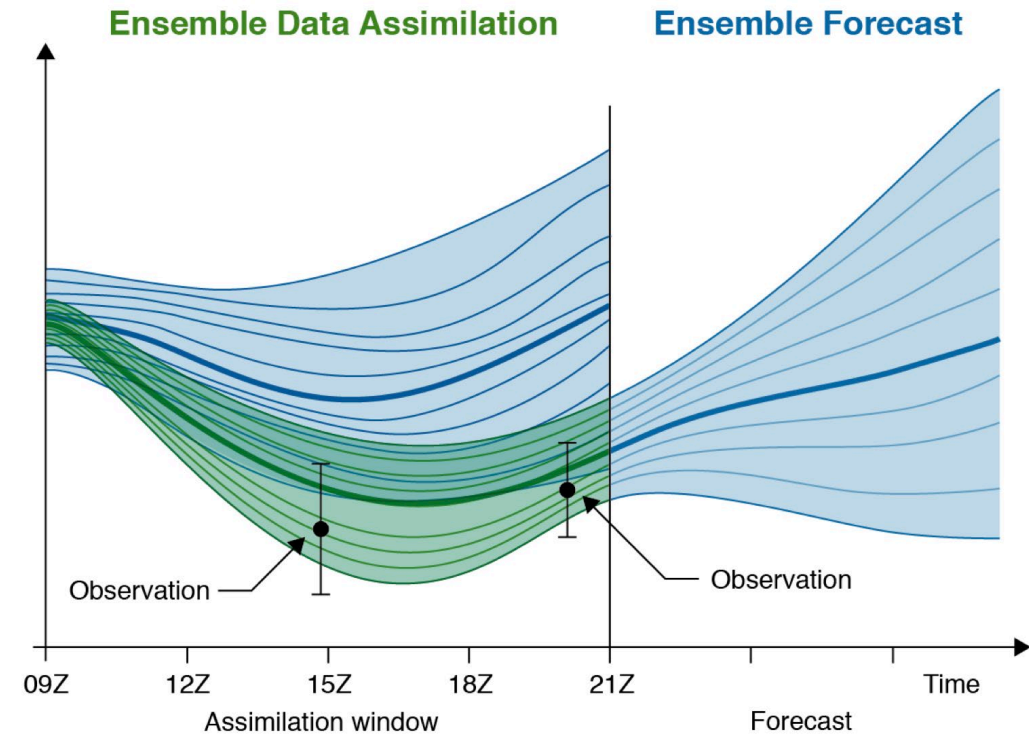
**Ensemble DA** (EDA) produces a collection of many perturbed MAP optimisations

$$\mathbf{x}^{a,(i)} = \arg \max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}^{(i)}, \mathbf{x}^{b,(i)}, \boldsymbol{\eta}^{(i)}), \quad i = 1, \dots, N.$$

**$N = 50$  at ECMWF**

As well as providing ICs for ensemble forecasting, the EDA provides **flow-dependent estimates of analysis and background error uncertainty** (Isaksen *et al* (2010)),

$$\mathbf{B}(t) \approx \frac{1}{N-1} \sum_{i=1}^N [\mathbf{x}^{b,(i)}(t) - \bar{\mathbf{x}}^b(t)] [\mathbf{x}^{b,(i)}(t) - \bar{\mathbf{x}}^b(t)]^T.$$



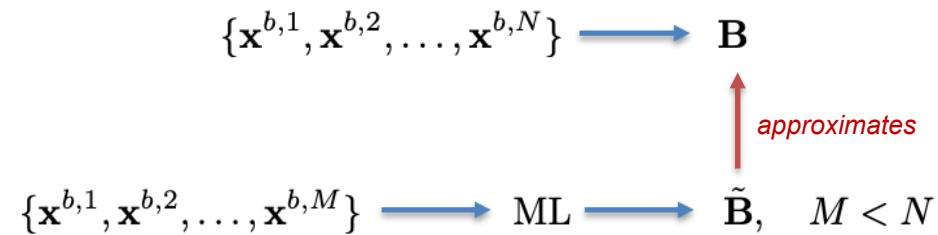
# Hybrid 4DVar-ML EDA

**Aim** – Develop an EDA emulator and integrate it into the operational EDA workflow.

Two complementary directions:

## 1. EDA statistics

- Emulate full-EDA statistics using fewer members.



## 2. EDA members generation

- Generate EDA members that effectively mimic the full 4DVar EDA.

- Full 4DVar EDA ensemble  $\{\mathbf{x}^{a,1}, \mathbf{x}^{a,2}, \dots, \mathbf{x}^{a,N}\}$
- Hybrid 4DVar and ML EDA ensemble  $\{\mathbf{x}^{a,1}, \mathbf{x}^{a,2}, \dots, \mathbf{x}^{a,M}\} \cup \{\tilde{\mathbf{x}}^{a,j_1}, \tilde{\mathbf{x}}^{a,j_2}, \dots, \tilde{\mathbf{x}}^{a,j_{N-M}}\}$

ML emulated members

# Hybrid 4DVar-ML EDA

## Current status

### 1. EDA statistics

- Developed and evaluated a prototype emulator for operational 50-EDA sampled stDev  $\Sigma(t)$  (**es**), using input ensemble size of **n=5** (*ECMWF Tech Memo 936*).



### 2. EDA members generation

- Developed and validated a ML methodology for EDA analysis perturbations using the Lorenz-96 low-order model as a testbed.
- A scaled prototype of the method at higher resolution is under development (using Anemoi) and testing.
  - Currently working on integrating the model into our EDA workflow for testing, evaluation and improvements.

## EDA statistics

# EDA variance emulator

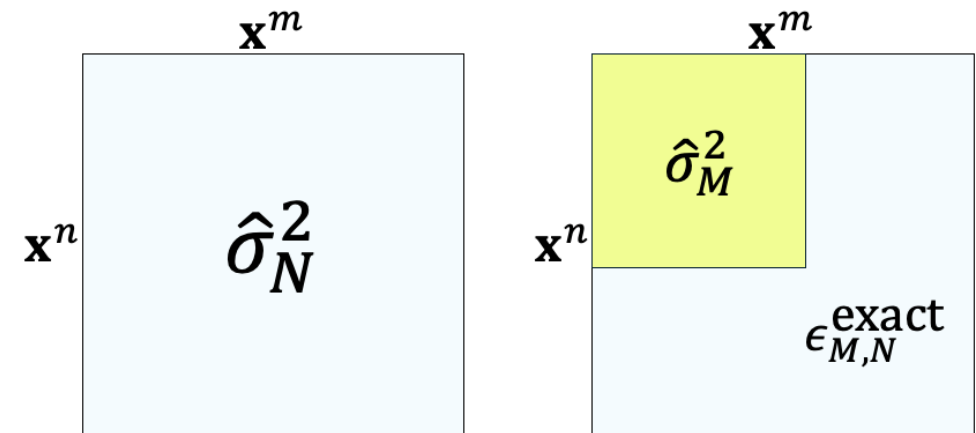
See *ECMWF Tech Memo 936* for details; doi: [10.21957/0b7e4d4426](https://doi.org/10.21957/0b7e4d4426)

- The variance  $\Sigma$  part of  $\mathbf{B}$  is just the pointwise unbiased sample variance

$$\hat{\sigma}_N^2 := \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}^{b,(n)} - \bar{\mathbf{x}}^b)^{\otimes 2} = \frac{1}{2N(N-1)} \sum_{n,m=1}^N (\mathbf{x}^{b,(n)} - \mathbf{x}^{b,(m)})^{\otimes 2}$$

- Fix a small ensemble estimate,  $M < N$

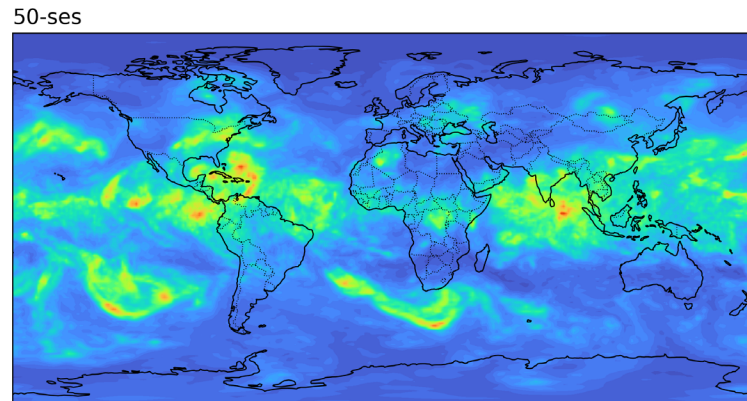
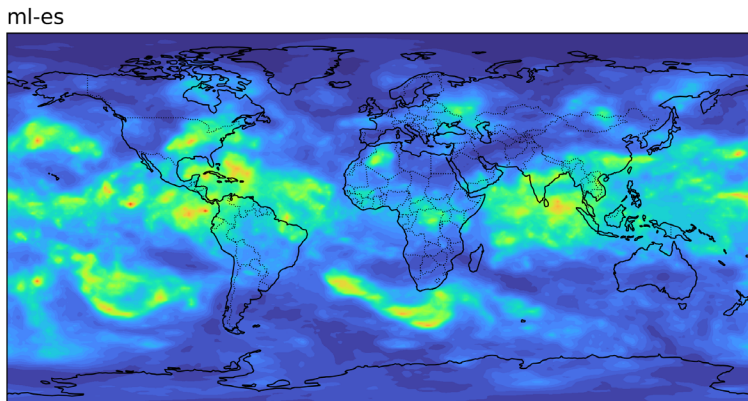
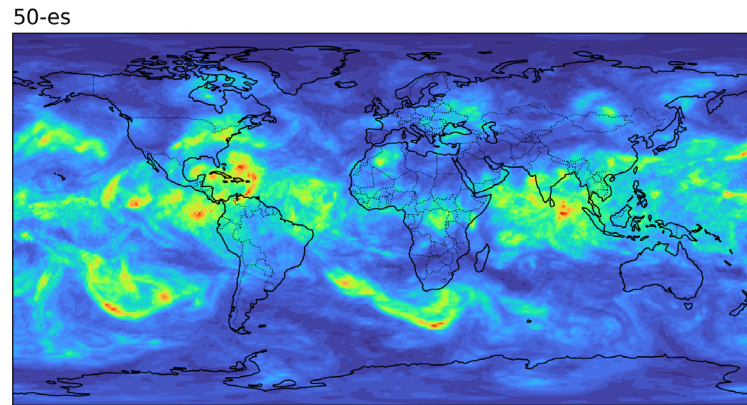
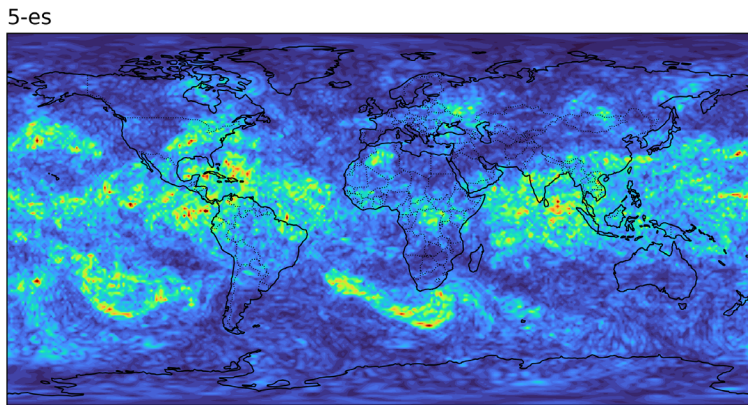
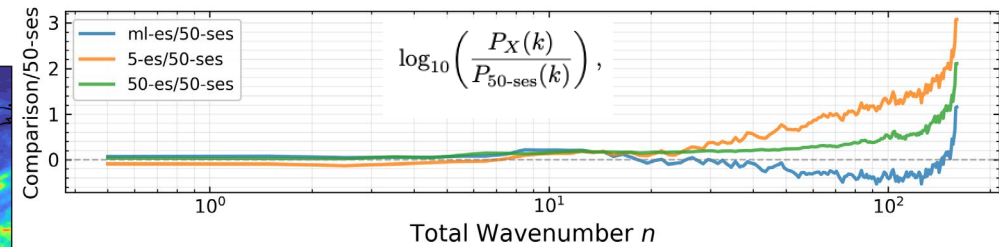
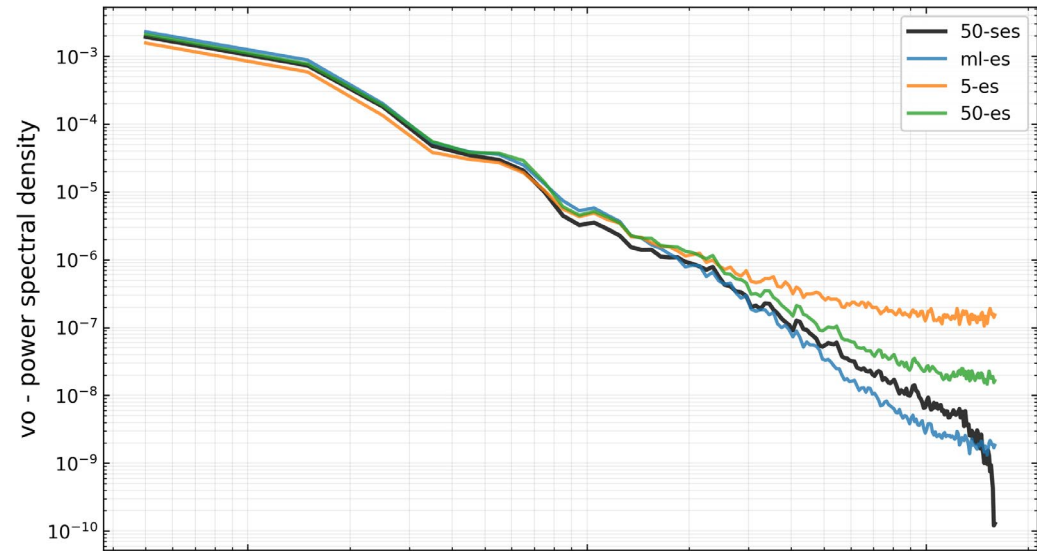
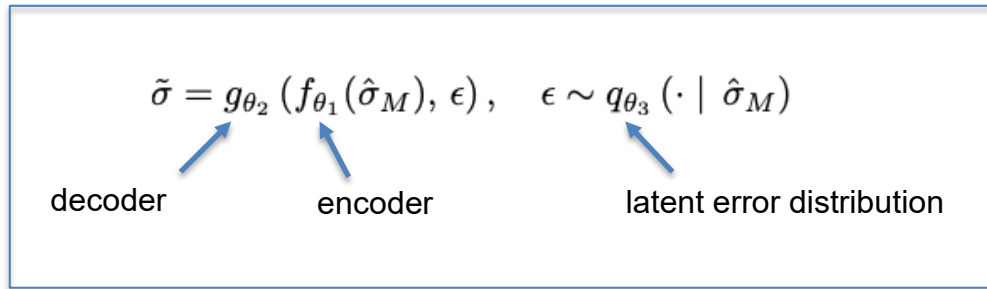
where



- The proposed model takes as input  $\hat{\sigma}_M$ . Its output is a probability distribution on  $\hat{\sigma}_N$

$$\hat{\sigma}_M \mapsto \mathbb{P}_\theta(\sigma_N \mid \hat{\sigma}_M).$$

# EDA variance emulator



**M=5, N=50**  
 N80 grid, 137 model levels  
 6 variables: (vo, div, t, Insp, q, o3)

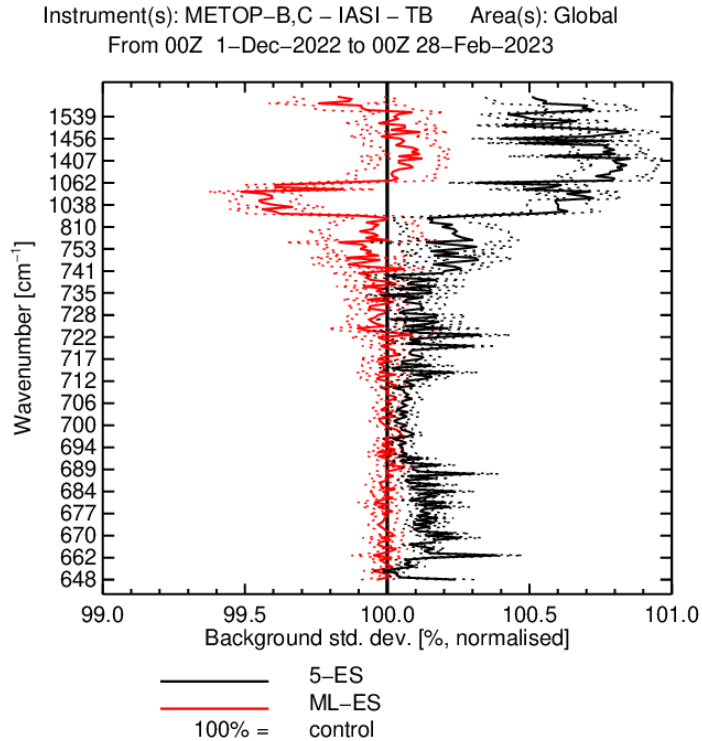
The training dataset spans  
 01/01/2023 – 31/12/2023.

Each model has  
 < 200,000 learnable parameters

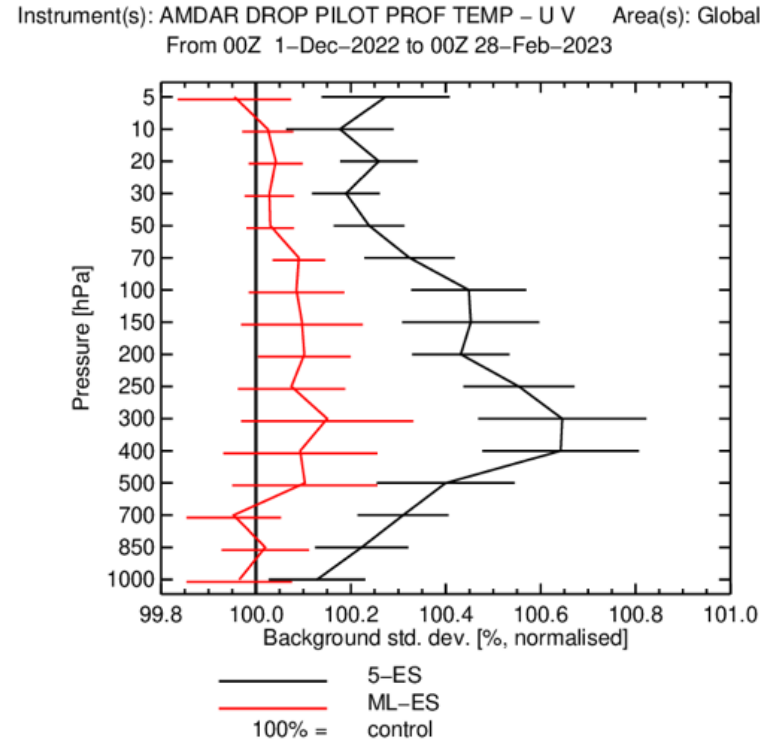
e.g. Vorticity es, ml74 (~200hPa)  
 21:00 UTC 01 June 2022

# Normalised observation statistics (obstats)

Satellite instruments



Conventional instruments



- Perfect replication is not possible due to information theoretic lower bounds

$$\mathbb{E}[(\bar{\sigma} - \hat{\sigma}_N)_{i,(\lambda,\phi)}^2] \geq \mathbb{E}[(\mathbb{E}(\hat{\sigma}_N | \mathcal{F}_M) - \hat{\sigma}_N)_{i,(\lambda,\phi)}^2],$$

any estimator from M samples

Optimal with respect to mean square error

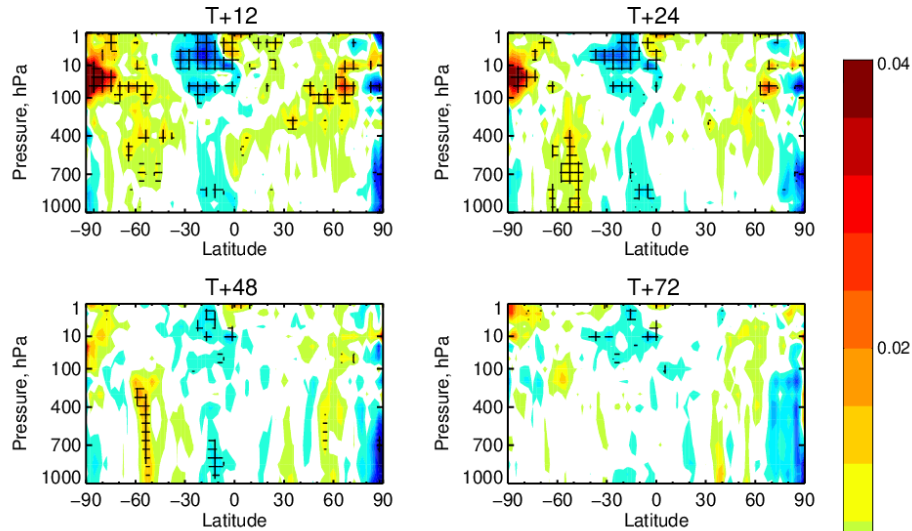
- 4DVar optimisation is **continuous** with respect to **B**, i.e. a small change in **B** means a small change in the resulting analysis increment.

# Normalised change in RMS forecast error

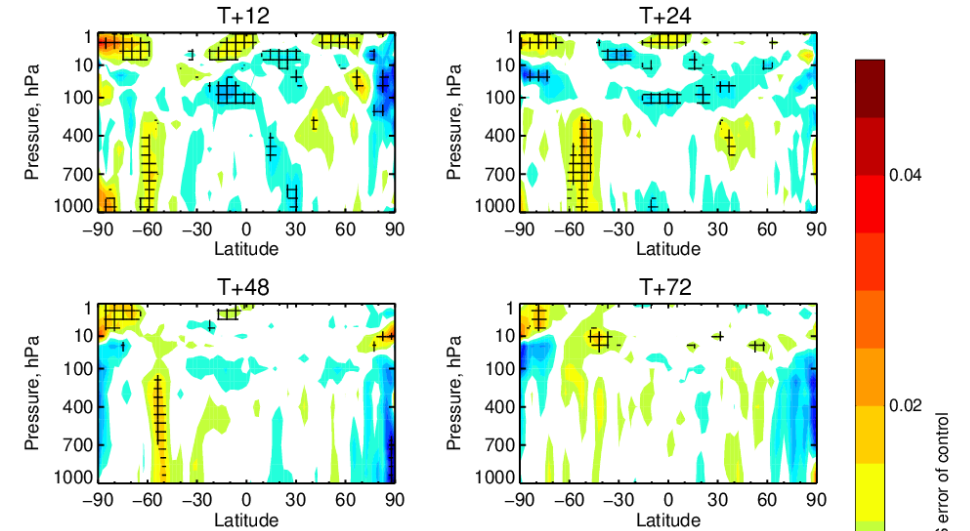
ML-es  
from 5  
members

vs 50-es

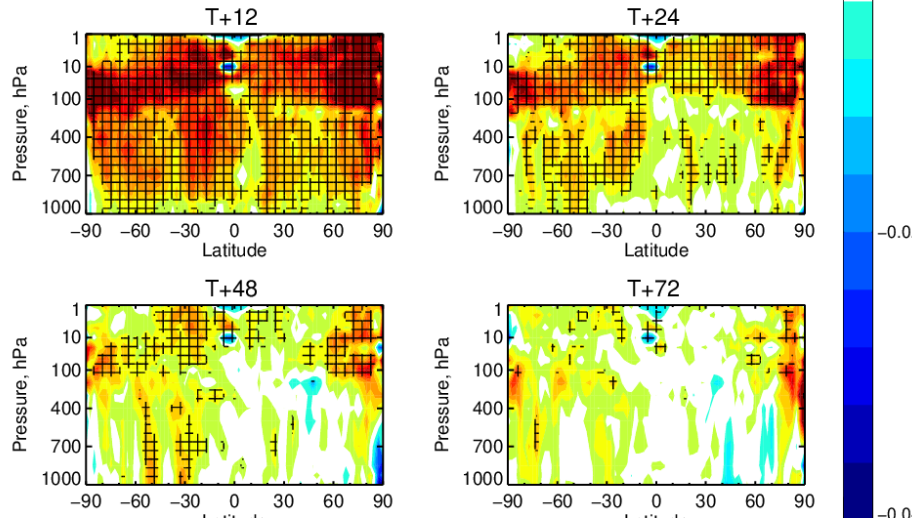
Change in RMS error in VW (ML-BAL-ES-5-control)  
1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.  
Cross-hatching indicates 95% confidence with Sidak correction for 20 independent tests.



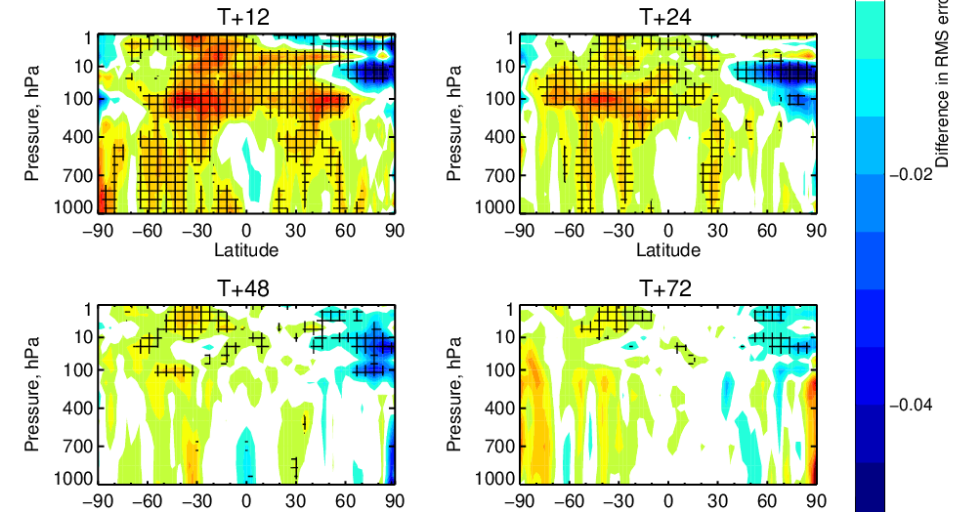
Change in RMS error in Z (ML-BAL-ES-5-control)  
1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.  
Cross-hatching indicates 95% confidence with Sidak correction for 20 independent tests.



Change in RMS error in VW (ES-5-control)  
1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.  
Cross-hatching indicates 95% confidence with Sidak correction for 20 independent tests.



Change in RMS error in Z (ES-5-control)  
1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.  
Cross-hatching indicates 95% confidence with Sidak correction for 20 independent tests.



5-es  
vs 50-es



## EDA members generation

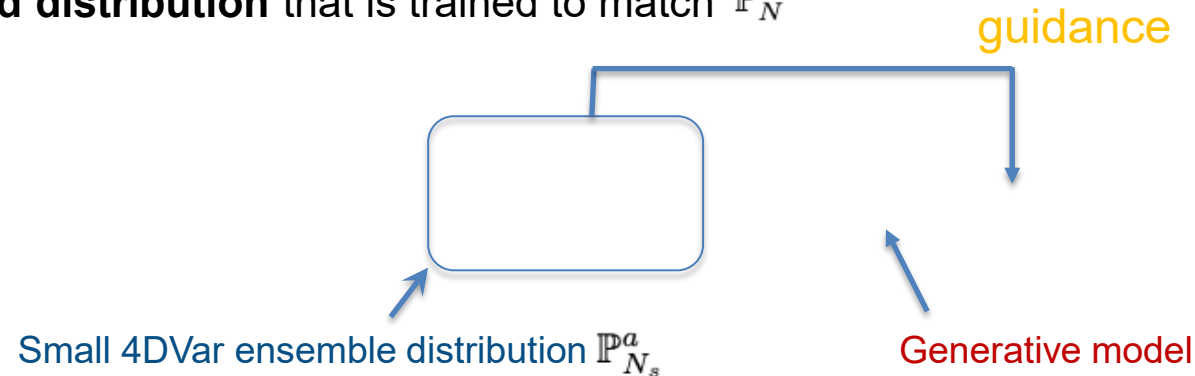
# EDA Emulation

The EDA system provides empirical background and analysis distributions

$$\mathbb{P}_N^b := \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}^{b,(i)}}, \quad \mathbb{P}_N^a := \frac{1}{N} \sum_{i=1}^N \delta_{\phi(\mathbf{x}^{b,(i)} | \mathbf{y}_{\text{obs}})}$$

where  $\phi$  denotes the 4DVar transformation applied to the background states.

We wish to learn a **hybrid distribution** that is trained to match  $\mathbb{P}_N^a$

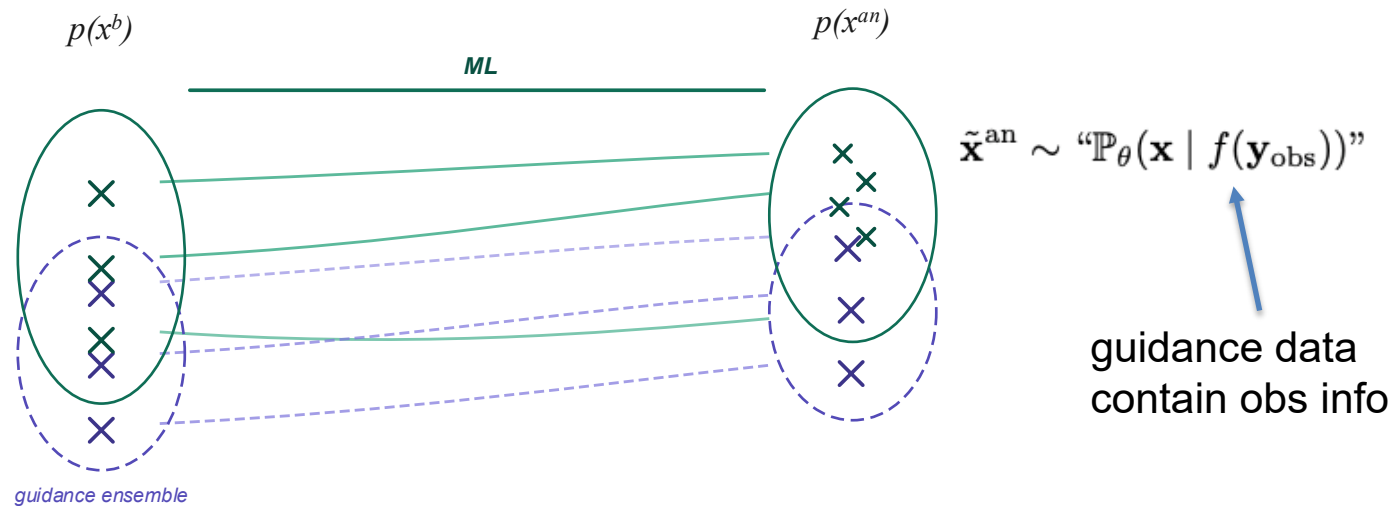


We propose a **stochastic flow** model (stochastic transport). Stochasticity enables cheap ensemble enlargement.

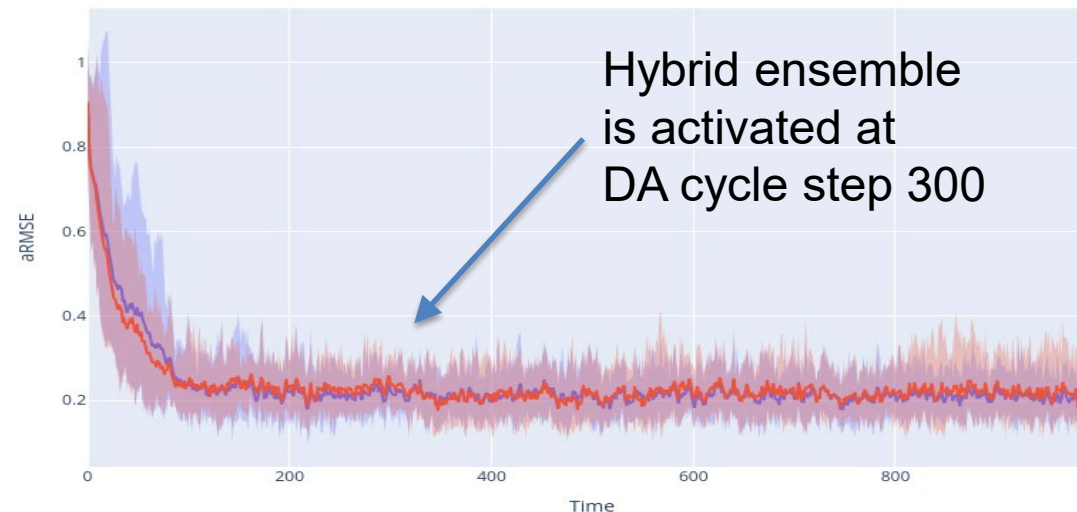
$$\phi_t^{(j)} = u_\theta(\phi_t^{(j)}, \mathbb{P}_t, t | \mathbb{P}_{N_s}^a) dt + \sigma_\theta(\phi_t^{(j)}, \mathbb{P}_t, t | \mathbb{P}_{N_s}^a) \circ dW_t, \quad \phi_0^{(j)} \in \left\{ \mathbf{x}^{b,(k_j)} \right\}_{j=1}^{N-N_s}$$

# EDA Emulation

Intuitive picture – the small ensemble of 4DVar analyses guides the generative model, so that observations are assimilated indirectly.



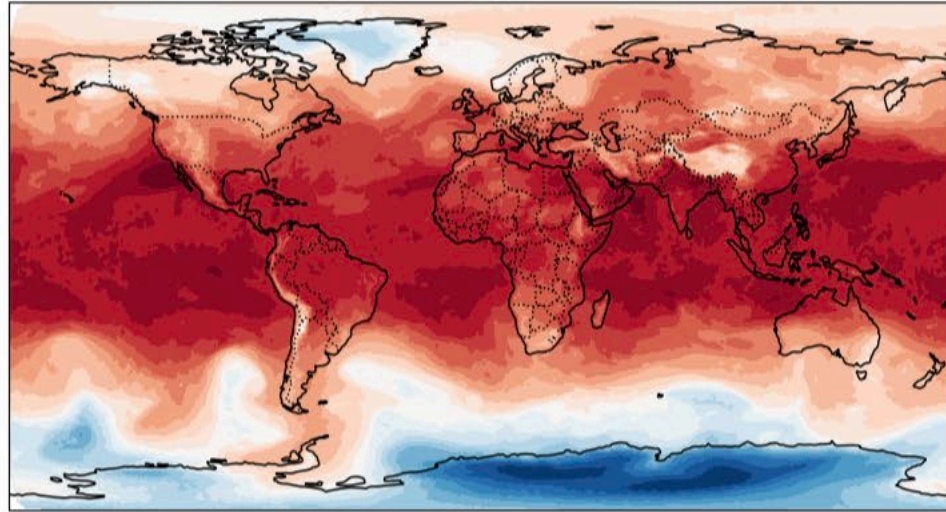
- **Lorenz96 ETKF** (from Alban Farchi)
  - dim=40,  $N=20$  (theoretical minimum ens. size 14)
  - Inflation only, no localisation
- Hybrid ensemble  $N_s = 10$  and 10 emulated.
- Neutral RMSE scores – **Hybrid EDA** vs **Full EDA**



# EDA Emulation: First Results

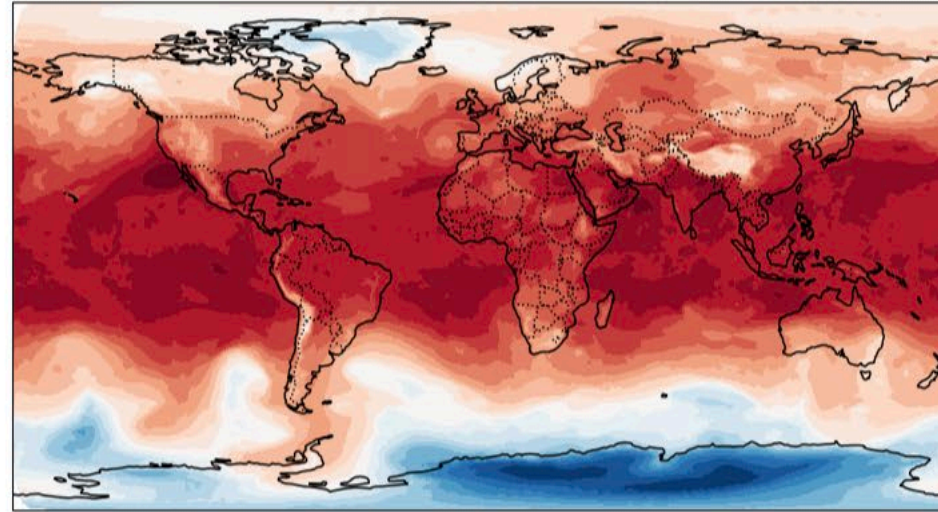
$a_n[i-1] \rightarrow a_n[i]$

ML Prediction - t\_100



emulated member

Target - t\_100



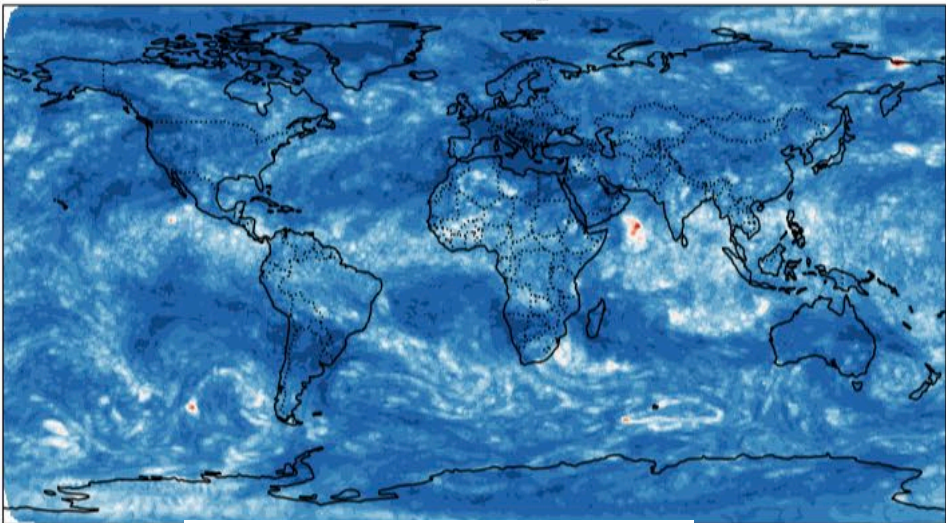
EDA member

Model trained using 4 ensemble members

ECMWF OD Dataset  
2024.07.01 – 2025.06.30

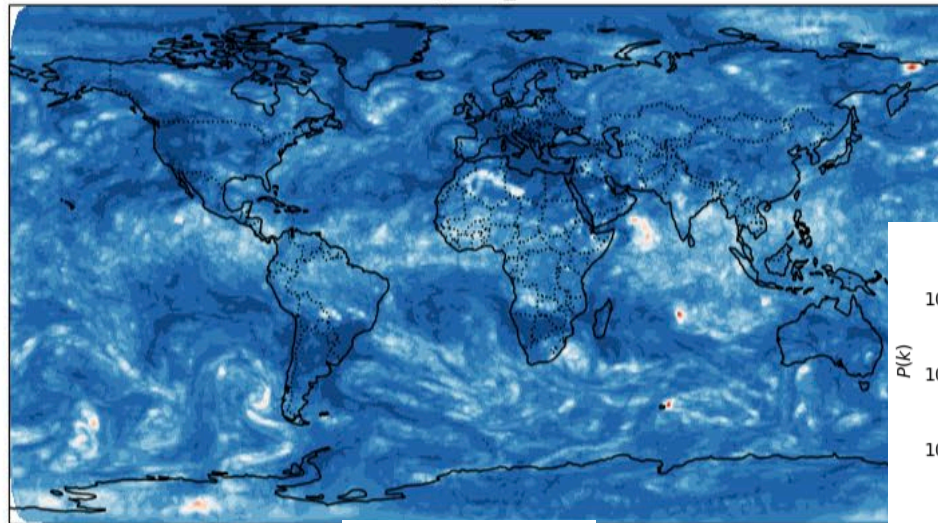
(t, vo, div, q, o3, Insp) full model levels, N128 grid.

ML Prediction - t\_100



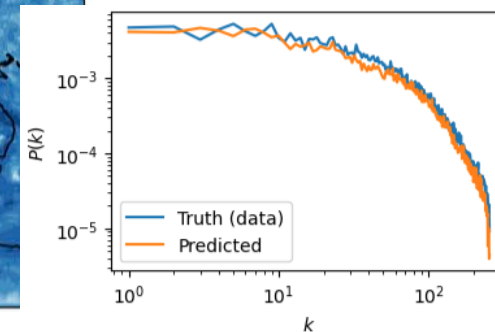
es (14 EDA + 36 emulated)

Target - t\_100



es (50 EDA)

Resolution is consistent with 50R1 increments.



Thank you for your attention.